

DISCRETE INDEPENDENT COMPONENT ANALYSIS (DICA) WITH BELIEF PROPAGATION

Francesco A. N. Palmieri and Amedeo Buonanno

Dipartimento di Ingegneria Industriale e della Informazione,
Seconda Università degli Studi di Napoli (SUN), Aversa, Italy
{francesco.palmieri; amedeo.buonanno}@unina2.it

ABSTRACT

We apply belief propagation to a Bayesian bipartite graph composed of discrete independent hidden variables and discrete visible variables. The network is the Discrete counterpart of Independent Component Analysis (DICA) and it is manipulated in a factor graph form for inference and learning. A full set of simulations is reported for character images from the MNIST dataset. The results show that the factorial code implemented by the sources contributes to build a good generative model for the data that can be used in various inference modes.

Index Terms— Bayesian Networks; Belief Propagation; ICA;

1. INTRODUCTION

Bi-directional information flow in belief propagation networks is becoming a very popular framework in many signal processing applications [1][2] because inference and learning can be easily manipulated with a small set of rules. Generally Bayesian models aim at capturing the hidden structure that may underly observed data through the assumption of a network of random variables that are only partially, or occasionally, visible [3].

Independent Component Analysis (ICA) is a popular signal processing framework in which observed data are mapped to, or generated from, independent hidden sources variables [4]. The variables are typically continuous and the transformation between sources and visible variables is linear. ICA has been used in many applications for signal separation and for analyzing signals and images [4]. ICA filters, trained on real images, seem to converge to patterns that resemble the receptive fields found in the neural visual cortex [5].

In this paper we explore the possibility of using the generative model of the ICA on discrete variables. The Bayesian model is constrained to a finite number of discrete hidden

sources (factorial code) that feed the visible variables, also discrete. Even if there are computational difficulties that naturally emerge in dealing with the product space of discrete alphabets, we find that even limiting our attention to tractable small sizes, the DICA framework clearly shows some potential in the applications, perhaps as a building block of more complex architectures. Discrete Component Analysis (DCA) has also been discussed by Buntine et al. [7] with reference to different models.

We reduce the DICA architecture to a Bayesian factor graph in the so-called reduced normal form (see [9] and reference therein) that includes only simple interconnected blocks. We experiment with belief propagation on this architecture using images extracted from the MNIST dataset [12]. We show that the DICA network nicely converges after learning to a generative model that reproduces accurately the image set.

In Section 2 the Bayesian model is presented and in Section 3 its discrete version is transformed into a factor graph for belief propagation. The various modes of inference are discussed in Section 5 and learning in Section 6. The simulations for unsupervised mapping of the MNIST images are reported in Section 6 with the addition of the label variable in Section 7. The conclusions are in Sections 8.

2. THE BAYESIAN MODEL

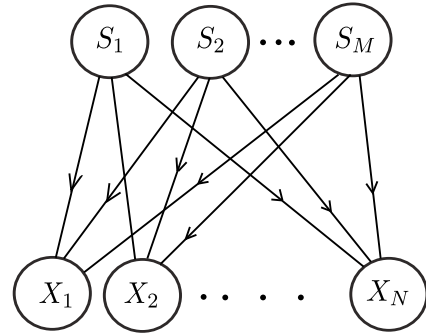


Fig. 1. The Bayesian Graph for M independent sources

Work partially supported by PON03PE-00185-1 - C3ISR, with Ministero dell'Istruzione dell'Università e della Ricerca; and PON03PE-00185-2 - MAR.TE., with Consorzio Nazionale Interuniversitario per le Telecomunicazioni (CNIT).

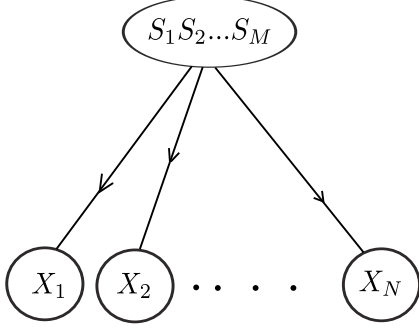


Fig. 2. The Bayesian Graph for M independent sources after the sources have been grouped (married).

In this paper we focus on the generative model depicted as the bi-partite graph of Figure 1 with M independent source variables S_1, S_2, \dots, S_M (hidden). The main variables X_1, X_2, \dots, X_N (visible), are connected to the source variables via the factorization

$$p(X_1 X_2 \dots X_N S_1 S_2 \dots S_M) = p(X_1 | S_1 S_2 \dots S_M) p(X_2 | S_1 S_2 \dots S_M) \dots p(X_N | S_1 S_2 \dots S_M) p(S_1) p(S_2) \dots p(S_M) \quad (1)$$

Note that X_1, X_2, \dots, X_N to be conditionally independent, must be conditioned on the whole set of sources, even if their marginal distribution factorizes: $p(S_1 S_2 \dots S_M) = p(S_1) p(S_2) \dots p(S_M)$. This appears to be the most general model for independent hidden sources that underly a set of dependent variables X_1, X_2, \dots, X_N . When $M = 1$, the system degenerates into a single-variable latent model [2].

One way of solving for the probability functions involved in the Bayesian model is to group (marry) the source variables (parents) [8] as in Figure 2. Note that the Bayesian graph does not show that the source variables are marginally independent. This is made more explicit in the factor graph representation that will follow.

2.1. Generative model for classical ICA

Independent Component Analysis is obtained when all the variables $x_1, x_2, \dots, x_N, s_1, s_2, \dots, s_M \in \mathcal{R}$ and the conditional probability density functions $p(x_i | s_1 s_2 \dots s_M)$ are constrained to depend on linear combinations of s_1, s_2, \dots, s_M . More specifically, the typical assumption is that the linear combinations contribute to the means of X_1, \dots, X_N and the dispersion around the mean is spherical and follows a Gaussian distribution

$$p(x_i | s_1 s_2 \dots s_M) = \mathcal{N}(x_i; \mathbf{a}_i^T \mathbf{s}, \sigma^2), \quad i = 1, \dots, N, \quad (2)$$

where the vector \mathbf{s} contains all the source values $\mathbf{s}^T = [s_1 s_2 \dots s_M]$ and \mathbf{a}_i is the i th column of the $N \times M$ coefficient matrix $A = [\mathbf{a}_1 \mathbf{a}_2 \dots \mathbf{a}_M]$ [5]. More compactly $p(\mathbf{x} | \mathbf{s}) = \mathcal{N}(\mathbf{x}; A^T \mathbf{s}, \sigma^2 I_N)$, where $\mathbf{x}^T = [x_1 x_2 \dots x_N]$. The

sources' pdfs $p(s_1), p(s_2), \dots, p(s_M)$ can follow various distributions that go from uniform to laplacian [5]. Typically for the model to be identifiable, the sources cannot be Gaussian (except perhaps for one out of M).

Unfortunately when ICA is used as a generative model it is hard to produce realistic images even when experimental densities are used as density sources [5]. Structured patches are easy to obtain, but they do not resemble the complex structures found in natural images. The reason is that independent continuous sources do not carry the necessary structure to assemble the ICA into the complex structures found in natural images. We report a simulation in the following that seems to confirm these results. Attempts have been made to use the ICA in two-layer architectures [5]. However, it is not clear how to properly include non linearities (without non linearities the whole system would still be linear) and investigations in this direction are still in progress.

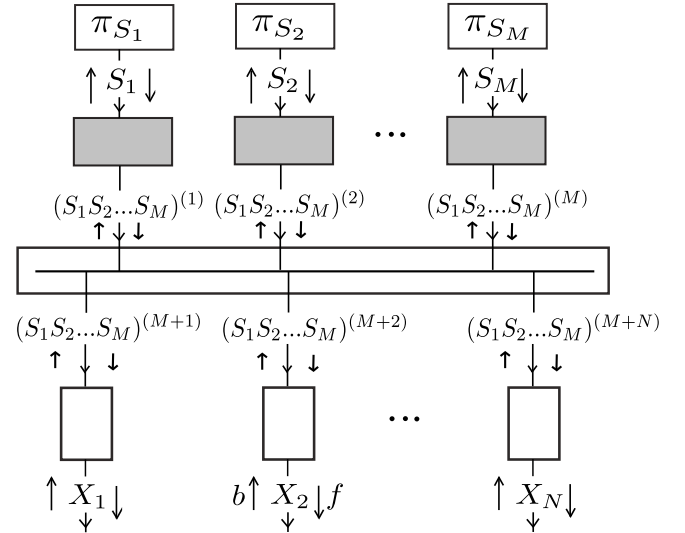


Fig. 3. The DICA model as a factor graph in reduced normal form. The shaded boxes represent the fixed matrices $P(S_1 S_2 \dots S_M | S_i)$, $i = 1, \dots, M$. The unshaded boxes represent the conditional probability matrices $P(X_j | S_1 S_2 \dots S_M)$, $j = 1, \dots, N$.

2.2. Discrete ICA

In this work we experiment on the unconstrained ICA model with discrete variables. More specifically we assume that both sources and visible variables take values in the finite discrete alphabets $\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_M, \mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_N$, with sizes $|\mathcal{S}_1|, |\mathcal{S}_2|, \dots, |\mathcal{S}_M|$ and $|\mathcal{X}_1|, |\mathcal{X}_2|, \dots, |\mathcal{X}_N|$.

The difficulties in dealing with such a model are clearly related to the computational complexity in the manipulation of the product space $\mathcal{S} = \mathcal{S}_1 \times \mathcal{S}_2 \times \dots \times \mathcal{S}_M$, that has size $|\mathcal{S}| = |\mathcal{S}_1| |\mathcal{S}_2| \dots |\mathcal{S}_M|$ (Figure 2). However, we find that even limiting our attention to small dimensionalities, i.e.

to few source variables and to small alphabets, the framework applied to natural images reveals quite interesting results. Furthermore, the basic architecture can be used as a building block for more complicated multi-layer Bayesian architectures (not discussed in this paper).

3. DICA IN REDUCED NORMAL FORM

Probability propagation and learning for the graph of Figure 1 can be handled in a very flexible way if we transform the model into a factor graph as in Figure 3. The graph is in the so-called *reduced normal* form (see [9] and references therein), that is composed only of one-to-one blocks, source blocks and diverters (these are equal constraint blocks that act like buses for belief propagation). One-to-one blocks are characterized by a conditional probability matrix and sources by a probability vector. We have often advocated the use of such a representation because it can be handled as a block diagram and it is amenable to distributed implementations. We have also designed a Simulink library for rapid prototyping [10].

More specifically for the DICA model, the source variables, that have prior distributions $\Pi_{S_1}, \dots, \Pi_{S_M}$, are mapped to the product space via the fixed row-stochastic matrices (shaded blocks)

$$\begin{aligned} P((S_1 S_2 \dots S_M)^{(1)} | S_1) &= \frac{|S_1|}{\prod_{i=1}^M |S_i|} I_{|S_1|} \otimes 1_{|S_2|}^T \otimes 1_{|S_3|}^T \otimes \dots \otimes 1_{|S_M|}^T, \\ P((S_1 S_2 \dots S_M)^{(2)} | S_2) &= \frac{|S_2|}{\prod_{i=1}^M |S_i|} 1_{|S_1|}^T \otimes I_{|S_2|} \otimes 1_{|S_3|}^T \otimes \dots \otimes 1_{|S_M|}^T, \\ &\dots \\ P((S_1 S_2 \dots S_M)^{(M)} | S_M) &= \frac{|S_M|}{\prod_{i=1}^M |S_i|} 1_{|S_1|}^T \otimes 1_{|S_2|}^T \otimes 1_{|S_3|}^T \otimes \dots \otimes I_{|S_M|}, \end{aligned} \quad (3)$$

where \otimes denotes the Kronecker product, 1_K is a K -dimensional column vector with all ones, and I_K is the $K \times K$ identity matrix. The conditional probability matrix is such that each variable contributes to the product space with its value and it is uniform on the components that compete to the other source variables. The blocks at the bottom of Figure 3 represent the $|S| \times |\mathcal{X}_j|$ conditional probability matrices $P(X_j | S_1 S_2 \dots S_M)$, $j = 1, \dots, N$, that with the source prior distributions are typically learned from data. Information flows in the network bi-directionally: for each branch variable there is a forward (f) and a backward (b) message, which are (or proportional to) discrete probability vectors. Messages are usually kept normalized for numerical stability. The variables connected to the diverter represent a replicated version of the same variable, but they all carry different forward and backward messages that are combined with the product rule [11]. Propagation through each one-to-one block follows the sum rule which in the variable direction is the matrix multiplication $f_{out} = P(out|in)^T f_{in}$ (already normalized) and in

the opposite direction $b'_{in} = P(out|in)b_{out}$ and $b_{in} = \frac{b'_{in}}{\sum b'_{in}}$ (normalization). After propagation for a number of steps equal to the graph diameter (if there are no loops), posterior probability p for a variable branch can be computed with the normalized product $p = \frac{f \odot b}{\sum (f \odot b)}$ (\odot denotes the element-by-element product of two vectors). For the reader not familiar with this framework, it should be emphasized that these simple rules are rigorous translation of marginalization and Bayes' theorem [11].

4. INFERENCE IN THE DICA GRAPH

The flexibility of this framework allows the use of the factor graph of Figure 3 in various inference modes. Information flow is bi-directional and assuming that all the parameters have been learned and that the unspecified messages are initialized to uniform distributions, we can use the DICA graph in:

(1) *Generation*: Source values are picked and are injected as forward delta distributions at S_1, S_2, \dots, S_M . After three steps of message propagation, the forward distributions are collected at the terminal variables X_1, X_2, \dots, X_N . They are the (soft) decoded version of the source values. Note that these are distributions that are typically displayed as their means or their argmaxes (see simulation results in the following).

(2) *Encoding*: Observed values for X_1, X_2, \dots, X_N are injected as delta backward distributions at the bottom. After three steps of message propagation, the backward distributions are multiplied with the forward at S_1, S_2, \dots, S_M . The normalized result is a (soft) *factorial code of the input*. The set of argmaxes of these distribution is the MAP decoding of the input.

(3) *Pattern completion*: Only a subset of values for X_1, X_2, \dots, X_N is available (there are erasures). The available values are injected at the bottom as delta backward distributions. For the missing values uniform densities are usually injected. After three steps of message propagation, forward distributions are collected at the bottom variables. For the observed variables the forward-backward products return just the deltas on the observations and provides no new information. At the unknown variables, the forward distribution is our best (soft) knowledge of that variable. Here too the means or the argmaxes can be used as a final result. The inference on the erasures is the synthesis of the information coming from the observations and the priors.

(4) *Error correction*: Available values for X_1, X_2, \dots, X_N may contain errors. They are presented as backward delta distributions at the bottom variables. After three steps of message propagation, forward distributions (or their means or argmaxes) are collected and used as corrections. No product with the backward is applied here because we do not know which component is reliable. In a similar scheme the values for X_1, X_2, \dots, X_N may be known softly via distributions that are injected at the bottom as backward messages.

Note that in both (3) and (4) also coded versions of the observations are available at the source branches.

5. LEARNING IN THE DICA GRAPH

To train the DICA system, we assume that a set of T examples is available for the visible variables $(x_1[n]x_2[n]\dots x_N[n]), n = 1, \dots, T$ (training set). Learning the system matrices for the bottom blocks and the vectors for the sources, is performed using an EM search. Various algorithms can be used, all inspired by a localized maximum likelihood cost function. The iterations are confined to each block and use only locally available forward and backward messages. Details on the learning algorithms for the factor graph in reduced normal form have been reported elsewhere and are omitted here for space reasons (see [6] [9] and references therein).

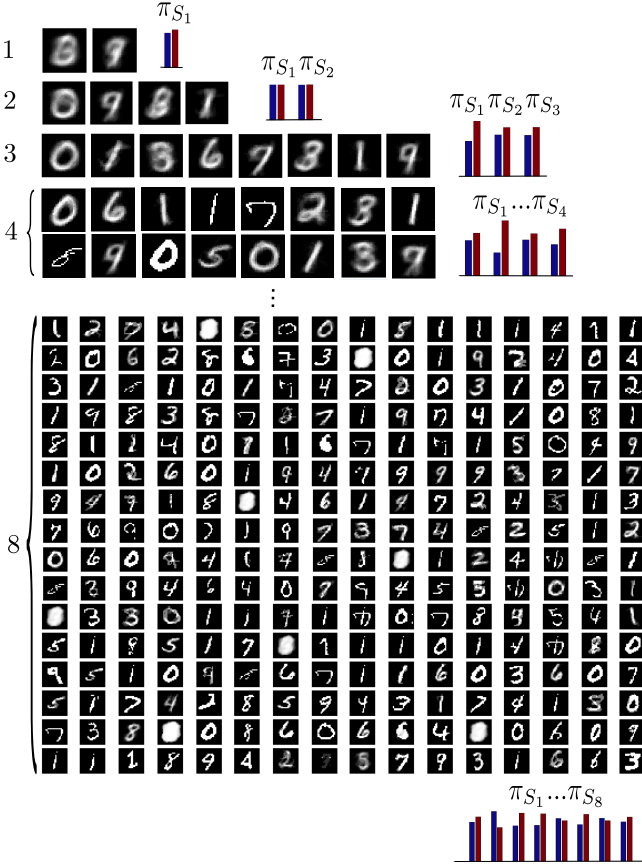


Fig. 4. Distribution means generated by the factorial code for increasing number of sources ($M = 1, 2, 3, 4, 8$). The bars show the learned source priors.

6. DICA SIMULATIONS

We report here a full set of simulations on the MNIST data set [12]. We have reduced the images to 28×28 binary pix-

els and extracted 500 images as our training set. In a first set of experiments we train the architecture of Figure 3 with all binary variables: $\mathcal{X}_j = \{x^0, x^1\}, j = 1, \dots, N$ ($N = 784$); $\mathcal{S}_i = \{s^0, s^1\}, i = 1, \dots, M$, for various number of sources $M = 1, 2, 3, 4, 8$. During learning the 500 images of the training set are presented as backward delta distributions on X_1, \dots, X_N , one time, with 5 cycles inside each block (the maximum likelihood algorithm inside each block is iterative [9]). Therefore for each order M we obtain the conditional probability matrices $P(X_j|S_1 \dots S_M), j = 1, \dots, N$, and the prior distributions $\pi_{S_1}, \dots, \pi_{S_M}$.

Generation: Figure 4 shows, for increasing M , the means of f_{X_1}, \dots, f_{X_N} when at the sources we inject the 2^M binary configurations in the forward messages f_{S_1}, \dots, f_{S_M} . Reported in the picture are also the learned priors. We note that, for larger number of sources, the product space (sizes 2,4,8,16,256), corresponds to increasingly accurate pattern memorization. For some characters, that are different in shape, the system builds separate representations. The source variables, independent by definition (factorial code), learn marginal distributions progressively less uniform as the number of sources increases (recall that the vector that represents $p(S_1, \dots, S_M)$ is the Kronecker product of the individual binary distributions and that even small non uniformities in the priors cause $p(S_1, \dots, S_M)$ to be highly non uniform).

Encoding: Figure 5 shows the typical results of presenting to the DICA graph of Figure 3, with $M = 8$, images from the test set (i.e. not included in the 500 images used for training) as backward delta distributions at X_1, \dots, X_N . In the third column the posterior distributions at the sources are shown (only the probability on the symbol s^1 is depicted). Here the DICA graph acts as an Encoder: the (soft) binary configurations are the factorial code of the presented images. Note that not all the codes are sharp. In the second column the mean of the forward distributions at X_1, \dots, X_N is also shown.

Decoding: In Figure 6 the same DICA graph is used as a soft decoder when smooth and sharp distributions are injected at the sources.

Pattern completion: Figure 7 shows the results of the same network when as backward at X_1, \dots, X_N we present images (from the test set) with 50 % of the pixels removed. For the erased pixels a backward uniform distribution is presented. The third and the fourth columns report the mean for the forward and the posterior distributions respectively. The network fills-in rather well the missing parts.

6.1. Continuous ICA on the same dataset

The natural question at this point is whether with continuous ICAs it would be possible to obtain similar results. The model is clearly very different, but on the same data set we have attempted a comparison. On the 500 MNIST images of the training set we have computed ICAs using the Fast ICA algorithm available for Matlab [13]. We have retained only

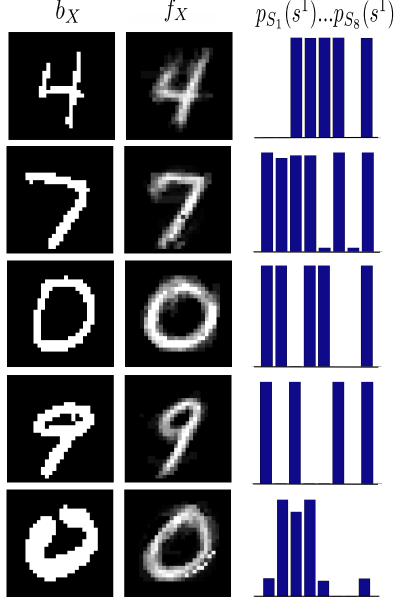


Fig. 5. Encoding of some images from the test set. Col. 1: images presented as delta backward distributions. Col. 2: means of the forward distributions. Col. 3: posterior probabilities at the sources (the bars represent $[p_{S_1}(s^1)...p_{S_8}(s^1)]$).

the first 8 components (largest variance) and estimated the output densities using average histograms. Random samples from these densities are used to generate the images through the inverse ICA [14]. Figure 8 shows the 8 masks and some generated images. The results confirm that, even if the ICA nicely represent bases for the data, with unconstrained independent samples at the sources, only average structures are generated. We have also tried with larger number of components and the obtained images look very similar. These results seem to be consistent with other experiments presented in the literature [14] for patches of natural images where only average textures are obtained. The linear ICA with independent unconstrained sources do not seem to be a generative model that preserves the structured composition of the training set.

7. DICA FOR CLASSIFICATION

The great flexibility of the factor graph framework allows to extend easily the architecture of the DICA graph to the one shown in Figure 9 where also a label variable C is included. The variable C belong to the finite alphabet $\mathcal{C} = c^0, c^1, \dots, c^9$ and it is attached directly, through a conditional probability matrix $P(C|S_1, \dots, S_M)$, to the product space diverter. Diverters in the reduced normal form act like probability pipelines [9].

Simulations have been performed on the same MNIST training set of 500 binarized images in the same mode as in the unsupervised experiments with the addition, during train-

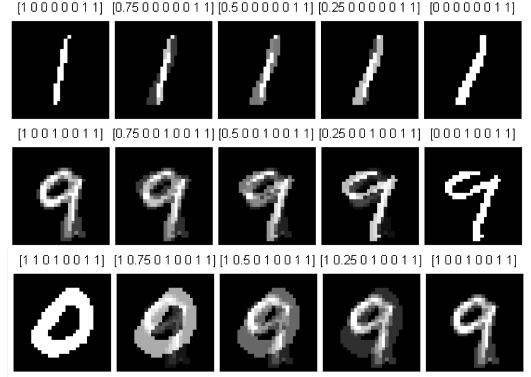


Fig. 6. Decoding for smooth forward distributions at the sources (in the brackets the probabilities $[f_{S_1}(s^1)...f_{S_8}(s^1)]$)

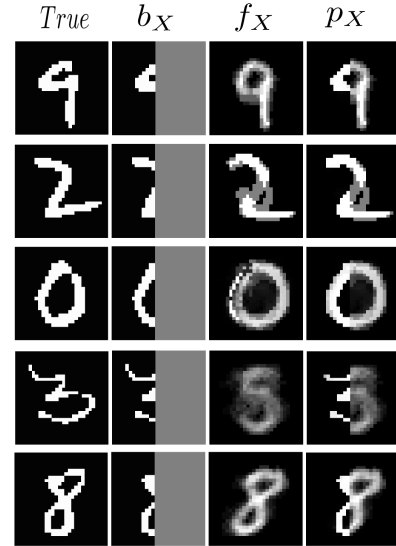


Fig. 7. Pattern completion of images from the test set after 50% removal.

ing, of the label information as a backward delta distribution. All the blocks, including now the probability matrix $P(C|S_1, \dots, S_M)$, are trained for $M = 8$. On the learned network, a typical recognition task on two images from the test set is shown in Figure 10. The bar graph represents simultaneously classification and encoding. Note how in the first row the network is naturally confused between c^4 and c^9 .

A generative experiment is also performed on this architecture with backward delta distributions injected at C . The results are shown in Figure 11. The images are the mean forward distributions at X_1, \dots, X_N and could be considered as the *prototypes* for the ten labels. The bar graphs are the corresponding simultaneous encoding at the sources.

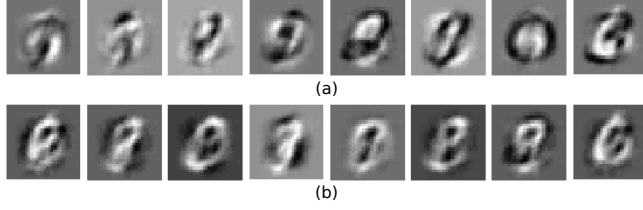


Fig. 8. Continuous ICA comparison: (a) 8 ICA masks for the Training Set (b) 8 generated images using the sources random values drawn from estimated output histograms.

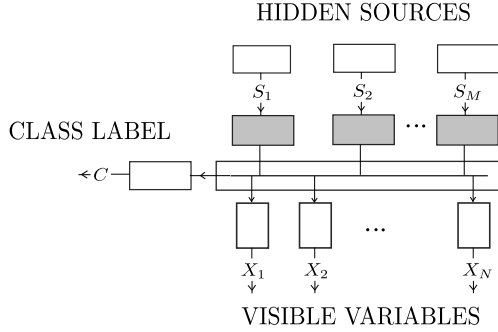


Fig. 9. The DICA model for classification

8. CONCLUSIONS

The simulations on the MNIST dataset with binary sources show that belief propagation in the DICA architecture, also with the addition of the label variable, provides a unified framework in which image data can be coded, generated and corrected in a very flexible way. We have also experimented on natural images on quantized patches obtaining very similar results, also when the sources have alphabet sizes greater than two. These results will be reported elsewhere. We are currently pursuing the use of this framework for building multi-layer architectures.

9. REFERENCES

- [1] M. I. Jordan, E. B. Sudderth, M. Wainwright, and A. S. Willsky, "Major advances and emerging developments of graphical models (and the whole special issue)," *IEEE Signal Processing Magazine*, vol. 17, November 2010, Special Issue.
- [2] D. Barber, *Bayesian Reasoning and Machine Learning*, Cambridge University Press, 2012.
- [3] J. Pearl, *Probabilistic reasoning in intelligent systems: networks of plausible inference*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1988.
- [4] A. Hyvarinen, J. Karhunen, and E. Oja, *Independent Component Analysis*, Wiley, New York, 2001.
- [5] Aapo Hyvriinen, Jarmo Hurri, and Patrick O. Hoyer, *Natural Image Statistics: A Probabilistic Approach to Early Computational Vision*, Springer Publishing Company, Incorporated, 1st edition, 2009.

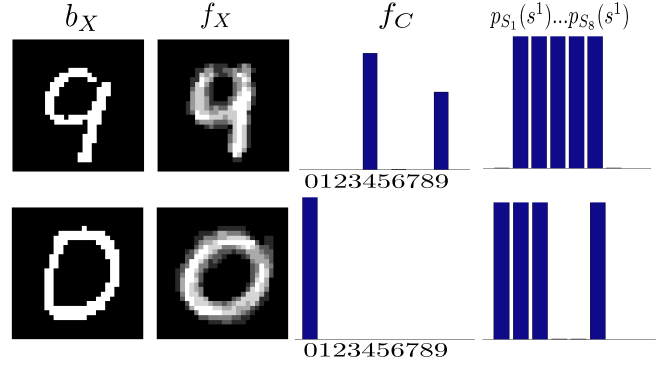


Fig. 10. Recognition task on two images from the test set

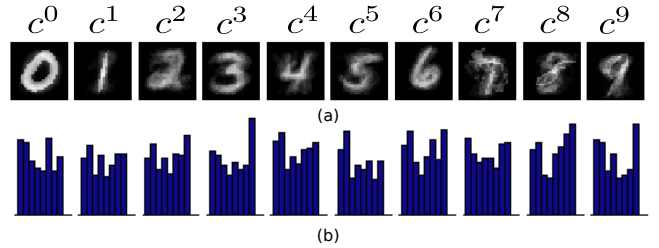


Fig. 11. Results of injecting backward deltas in C , $b_C(c) = \delta(c, c^i)$, $i = 0, \dots, 9$. (a) Means of the forward distributions (prototypes); (b) Posterior probabilities at the sources $[p_{S_1}(s^1) \dots p_{S_8}(s^1)]$ (encoding).

- [6] Francesco A. N. Palmieri, "Learning non linear functions with factor graphs," *IEEE Transactions on Signal Processing*, vol. 61, N. 7, pp. 4360–4371, 2013.
- [7] Wray Buntine and Aleks Jakulin, "Discrete component analysis," in *Subspace, Latent Structure and Feature Selection*, Craig Saunders, Marko Grobelnik, Steve Gunn, and John Shawe-Taylor, Eds., vol. 3940 of *Lecture Notes in Computer Science*, pp. 1–33. Springer Berlin Heidelberg, 2006.
- [8] S. L. Lauritzen, *Graphical Models*, Oxford, 1996.
- [9] F. A. N. Palmieri, "A comparison of algorithms for learning hidden variables in normal graphs," 2013.
- [10] A. Buonanno and F. A. N. Palmieri, "Simulink implementation of belief propagation in normal factor graphs," in *Proceedings of 2014 Workshop on Neural Networks, Vietri s.m.*, 2014.
- [11] H. A. Loeliger, "An introduction to factor graphs," *IEEE Signal Processing Magazine*, vol. 21, no. 1, pp. 28 – 41, jan. 2004.
- [12] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov 1998, <http://yann.lecun.com/exdb/mnist/>.
- [13] Hugo Gvert, Jarmo Hurri, Jaakko Srel, and Aapo Hyvriinen, "The fastica package for matlab," Available from: <http://research.ics.aalto.fi/ica/fastica/>.
- [14] Aapo Hyvriinen, "Statistical models of natural images and cortical visual representation," *Topics in Cognitive Science*, vol. 2, no. 2, pp. 251–264, 2010.